

# Prediction of PIK3CA mutation with gene expression

Jun Kang \*

\* Corresponding author: jkang.alien@gmail.com

## Introduction

A list

- Brevity
- Logic and clarity
- Clean typing

The problem

- PIK3CA mutation in selecting drug

Breast cancer with PIK3CA mutation has been approved to use PIK3CA inhibitor in hormone receptor positive HER2 negative subtype. [1] Prediction of PIK3CA mutation was done by gene expression data of TCGA.

- Predicting mutation by gene expression data
  - What for?
  - Feasible?
  - Cost?
- Varying frequency of PIK3CA mutation across cancer types
  - High in endometrial breast
  - not common other cancer type

Universal prediction using gene expression data across cancer type for certain kind of mutation. It's clinically useless now, but we want to explore the possibilities of the PIK3CA mutation prediction.

RNAseq can be widely used. The mutation status directly prediction

Previous study [2]

prediction of aberrant activation of a certain pathway vs a specific mutation The proposed solution

Prediction RAS pathway activation with gene expression data was done in previous study. [2] They trained pancancer The Cancer Genome Atlas (TCGA) data with a supervised elastic net penalized logistic regression classifier with stochastic gradient descent. The performance of their model was 84% with an area under the receiver operating characteristic (AUROC) curve and 63% with an area under the precision recall (AUPR) curve. We applied their modeling methods in prediction PIK3CA mutation.

# Materials and Methods

## Dataset

We used TCGA pancancer dataset. PIK3CA mutation data was get using cgdscr rpackage.[3] Gene expression data was get from GDAC firehose using RTCGAToolbox R package. [4] ER immunostain postive HER2 immunostain negative and/or SISH negative breast cancers are included. Data of invasive ductal carcinomas were used for training set and data of invasive lobular carcinoma were used for test set. Number of observations were 530 in training set and 188 in test set.

## Modeling

To narrow down potential predictors, Genes with a large the median absolute deviation (more than third-quartiles) were selected. 5000 out of 20502 genes were included in the modeling process.

Yeo-Johnson transformation was done to correct skewness. Centering and scaling were done. All preprocessing was done using recipe r package. [5]

Penalized logistic regression was applied to prediction modeling. 10-fold cross-validation with targe variable stratification was done over the hyperparameter grid:  $\lambda \{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10^0\}$ ,  $\alpha \{0.0, 0.25, 0.5, 0.75\}$ . Model performance was eevaluated with the area under receiver operating characteristic (ROC) and area under the ROC (AUROC) and precision recall (PR) and area under curve (AUPR). The modelin process was done with tidymodels rpackage. [6]

## Results

### selecting model

Cross-validation

The model showed best performance at  $\lambda = 0.01$  and  $\alpha = 1.0$  (Ridge regression).

### Prediction performance

The trainset AUROC was 0.93 and the testset AUROC was 0.84. The AUPR of trainset was 0.66 and the testset AUPR was 0.39. (Figure 1A)

### performance of each cancer type

Because tht prevalence of PIK3CA mutation is vary across the cancer type, the performance of each cancer type was investigated. The AUROC and AUPR were positively correlated between train set and test set in cancer type subanalysis. (Figure 1B) The AUPR was high in cancer type with high PIK3CA mutation rate such as colon, brest, Uterus cancer types. The AUROC did not correlated with PIK3CA mutation rate of each cancer types. (Figure 1C)

### Important predictors

## Discussion

Our model showed good performance to predict PIK3CA mutation of various cancer types.

Main message answers the question and main supporting evidence	70
Critical assessment opinions on - any shortcomings in study design	71
In this study, the split method was done for testset instead of independent testset.	72
Because there is no standard analysis method for the RNA-seq gene expression quantification, it is difficult to find gene expression data same as TCGA data.	73
Developing standard analysis method for gene expression data is necessary to apply this prediction model.	74
• limitations in methods	75
Case imbalance	76
• flaws in analysis	77
• validity of assumption	78
Comparison with other studies where inconsistencies are discussed	79
PIK3CA prediction model performed better than RAS activation prediction model in terms of AUROC. AUPR and the PIK3CA mutation rate interpretation. [7]	80
Conclusions comments on possible biological or clinical implications and suggestions for further research.	81
Evaluate the results - not the authors	82

## Figure legends

Figure 1. Summary of modeling results

(A) Left: receiver operating characteristic (ROC) curve right: precision recall (PR) curve of trainset and testset. The horizontal green line is the PIK3CA mutation rate (0.11) (B) Correlation between trainset and testset of area under receiver operating characteristic curve (AUROC) and area under precision recall curve (AUPR) among cancer type. The abbreviations are explained in S1 appendix. (C) Correlation between the PIK3CA mutation rate of area under receiver operating characteristic curve (AUROC) and area under precision recall curve (AUPR).

Figure 2. Coefficients of model

(A) Top 30 high coefficients of mRNA. (B) Coefficients of cancer types. The abbreviations are explained in S1 appendix.

## Supporting information

S1 Appendix.	99
S2 Figure.	100
S1 Table.	101
S2 Table.	102
S3 Table.	103

## References

1. André F, Ciruelos E, Rubovszky G, Campone M, Loibl S, Rugo HS, et al. Alpelisib for PIK3CA-Mutated, Hormone ReceptorPositive Advanced Breast Cancer. New	104
	107

England Journal of Medicine. Massachusetts Medical Society; 2019;380: 1929–1940.	108
doi:10.1056/NEJMoA1813904	109
2. Way GP, Sanchez-Vega F, La K, Armenia J, Chatila WK, Luna A, et al. Machine	110
Learning Detects Pan-cancer Ras Pathway Activation in The Cancer Genome Atlas.	111
Cell Reports. Elsevier; 2018;23: 172–180.e3. doi:10.1016/j.celrep.2018.03.046	112
3. Jacobsen A, Luna A. Cgdsr: R-based API for accessing the MSKCC cancer	113
genomics data server (CGDS). 2019.	114
4. Samur MK. RTCGAToolbox: A new tool for exporting TCGA Firehose data.	115
PLoS One. 2014;9(9):e106397.	116
5. Kuhn M, Wickham H. Recipes: Preprocessing tools to create design matrices.	117
2020.	118
6. Kuhn M, Wickham H. Tidymodels: Easily install and load the 'tidymodels'	119
packages. 2020.	120
7. Davis J, Goadrich M. The relationship between Precision-Recall and ROC curves.	121
Proceedings of the 23rd international conference on Machine learning. Pittsburgh,	122
Pennsylvania, USA: Association for Computing Machinery; 2006. pp. 233–240.	123
doi:10.1145/1143844.1143874	124